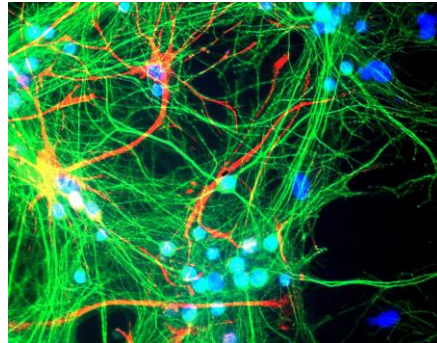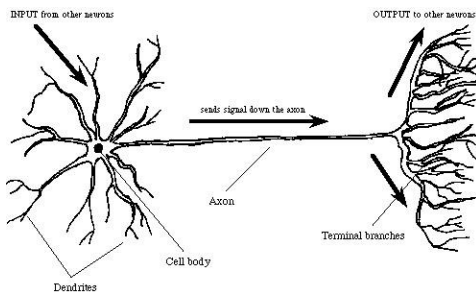# Neural Networks (part 1)

## Goals for this lecture

you should understand the following concepts
- perceptrons
- the perceptron training rule
- linear separability
- hidden units
- multilayer neural networks
- gradient descent
- stochastic (online) gradient descent
- activation functions
- sigmoid, hyperbolic tangent, ReLU
- objective (error, loss) functions
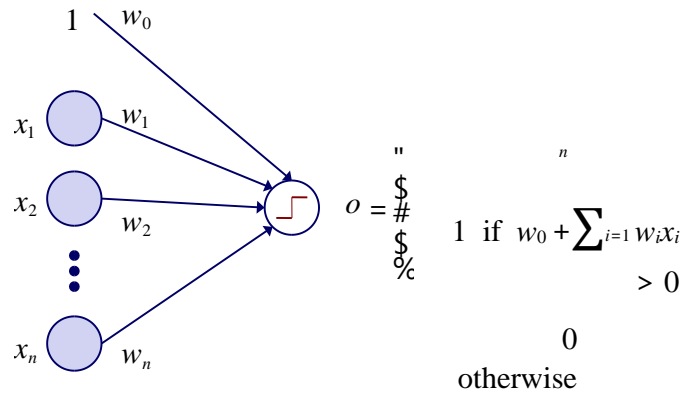- squared error, cross entropy

## Neural Networks

- a.k.a. *artificial neural networks*, *connectionist models*
- inspired by interconnected neurons in biological systems
- simple processing units
- each unit receives a number of real-valued inputs
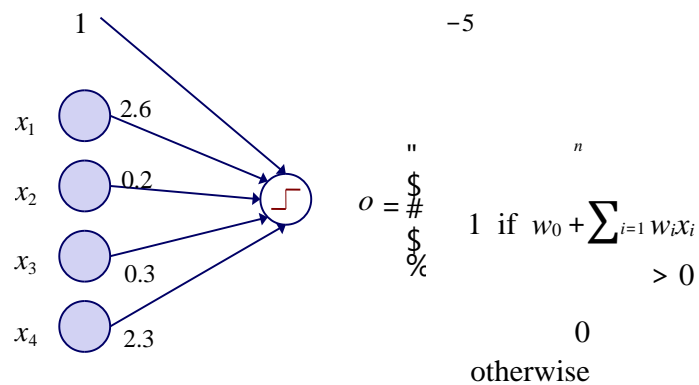- each unit produces a single real-valued output

# Perceptrons

$$o = \begin{cases} 1 & \text{if } w_0 + \sum_{i=1}^{n} w_i x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

*input units*:
represent given $x$

*output unit*:
represents binary classification

# Perceptron example

$$o = \begin{cases} 1 & \text{if } w_0 + \sum_{i=1}^{n} w_i x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

features, class labels are represented numerically

$$\mathbf{x} = \langle 1,0,0,1 \rangle \quad w_0 + \sum_{i=1}^{n} w_i x_i = -0.1 \quad o = 0$$

$$\mathbf{x}=1,0,1,1 \quad w_0 + \sum_{i=1} w_i x_i = 0.2 \qquad o=1$$

# Training a perceptron

1. randomly initialize weights

2. iterate through training instances until convergence

2a. calculate the output for
the given instance

$$o = \begin{cases} 1 & \text{if } w_0 + \sum_{i=1}^{n} w_i x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

2b. update each weight

$$\Delta w_i = \eta(y-o)x_i$$

$\eta$ is *learning rate*;

set to value << 1

$$w_i \leftarrow w_i + \Delta w_i$$
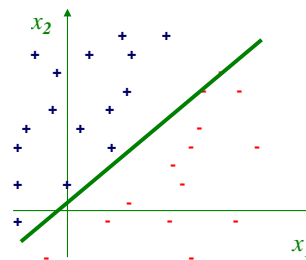
# Representational power of perceptrons

perceptrons can represent only *linearly separable* concepts

$$o = \begin{cases} 1 & \text{if } w_0 + \sum_{i=1}^{n} w_i x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$
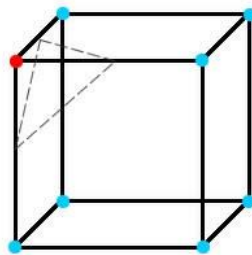
decision boundary given by: $1$ if $w_0$

$+ w_1 x_1 + w_2 x_2 > 0$  $w_1 x_1 + w_2 x_2 = -w_0$

$$x_2 = -\frac{w_1 x_1}{w_2} - \frac{w_0}{w_2}$$
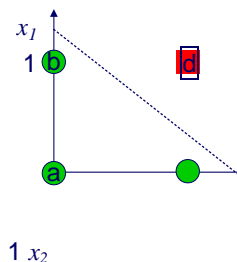


# Representational power of perceptrons

- in previous example, feature space was 2D so decision boundary was a line
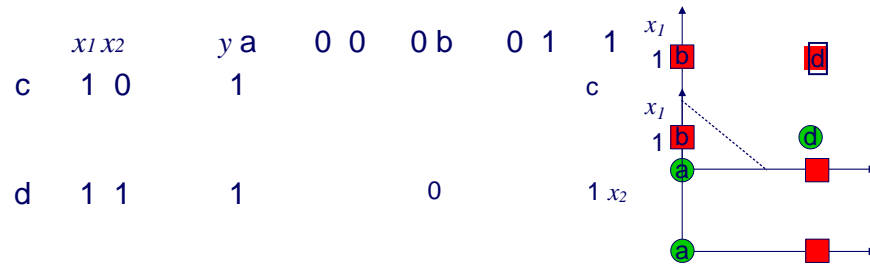- in higher dimensions, decision boundary is a hyperplane
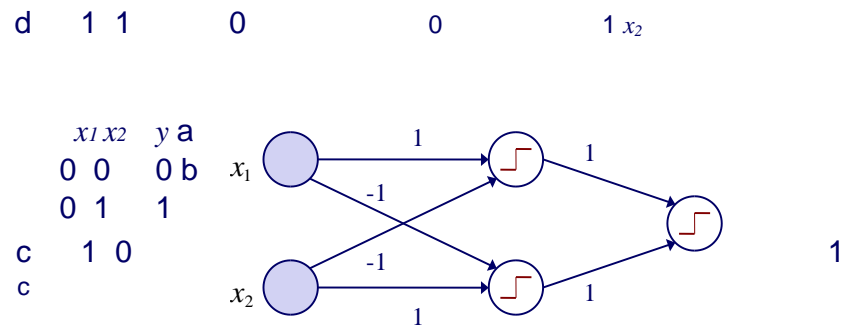


# Linearly separable functions

AND

| | $x_1$ | $x_2$ | $y$ |
|---|---|---|---|
| a | 0 | 0 | 0 |
| b | 0 | 1 | 0 |
| c | 1 | 0 | 0 |
| d | 1 | 1 | 1 |

|   | $x_1$ $x_2$ | $y$ a | 0 0 | 0 b | 0 1 | 1 | $x_1$ |
|---|---|---|---|---|---|---|---|
| c | 1 0 | 1 |   |   |   | c |   |

$x_1$

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| d | 1 1 | 1 | 0 | | 1 $x_2$ |

# Is XOR linearly separable?

| d | 1 1 | 0 | 0 | 1 $x_2$ |
|---|---|---|---|---|

| | $x_1$ $x_2$ | $y$ a |
|---|---|---|
| | 0 0 | 0 b |
| | 0 1 | 1 |
| c | 1 0 | 0 |
| c | | |

$x_1$  1  1

-1

-1  1

$x_2$  1

1

a multilayer perceptron
can represent XOR

assume $w_0 = 0$ for all nodes

# Example of multi-layer neural network



output units
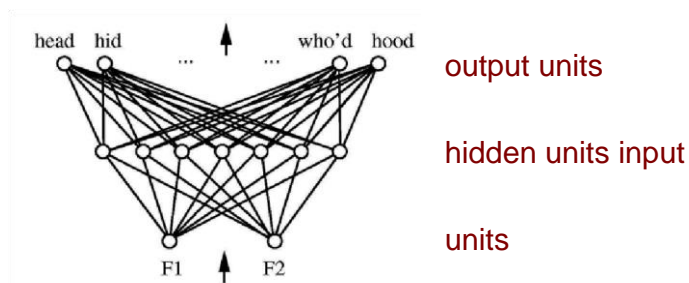
hidden units input

units

figure from Huang & Lippmann, *NIPS* 1988

input: two features from spectral analysis of a spoken sound output:

vowel sound occurring in the context "h__d"

# Decision regions of a multi-layer NN
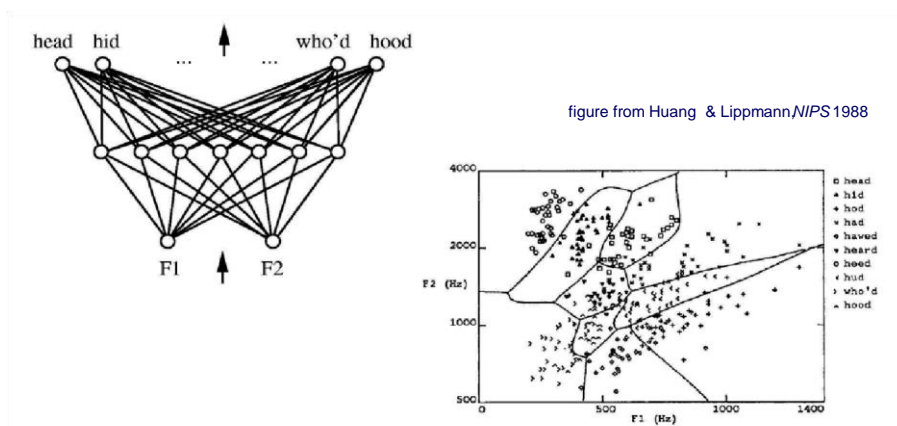
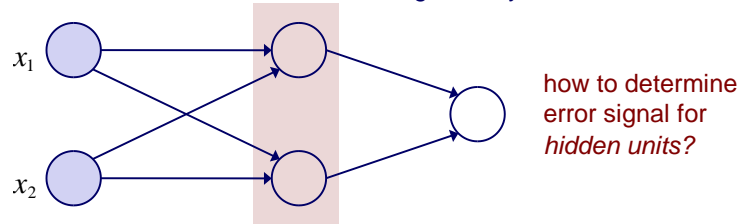

figure from Huang & Lippmann *NIPS* 1988

input: two features from spectral analysis of a spoken sound output:

vowel sound occurring in the context "h__d"

# Learning in a multi-layer NN

- work on neural nets fizzled in the 1960's
- single layer networks had representational limitations (linear separability)
- no effective methods for training multilayer networks

$x_1$

$x_2$

how to determine error signal for *hidden units?*

- revived again with the invention of *backpropagation* method [Rumelhart & McClelland, 1986; also Werbos, 1975]
- key insight: require neural network to be differentiable; use *gradient descent*
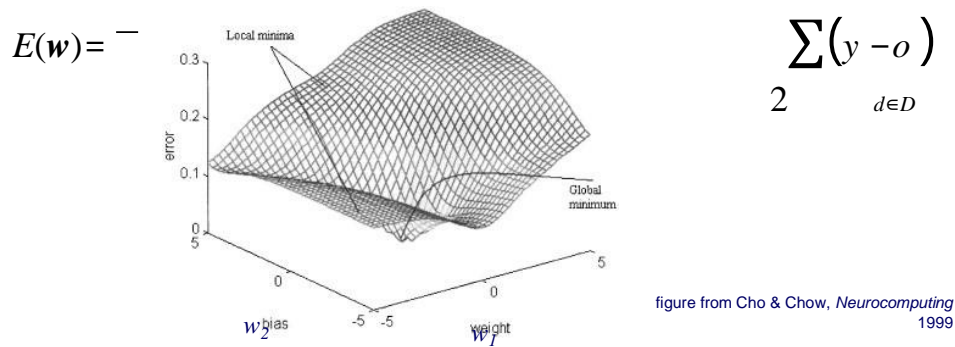
# Gradient descent in weight space

Given a training set                                              we can specify an

$$D=\left\{(\boldsymbol{x}^{(1)}, y^{(1)})\square(\boldsymbol{x}^{(m)}, y^{(m)})\right\}$$  error measure that is a function of our weight

vector      $\boldsymbol{w}$

$$1 \qquad {}^{(d)} \qquad {}^{(d)\,2}$$

$$E(\mathbf{w}) = \overline{\phantom{x}} \qquad \sum_{d \in D} (y - o)$$
$$2$$


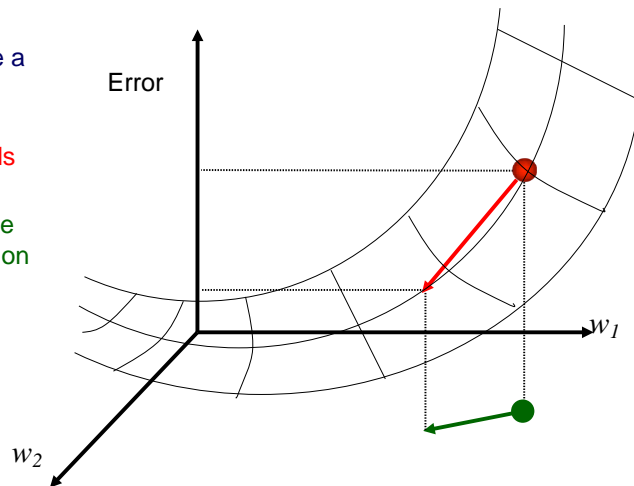
figure from Cho & Chow, *Neurocomputing* 1999

This *objective function* defines a surface over the model (i.e. weight) space

# Gradient descent in weight space

gradient descent is an iterative process aimed at finding a minimum in the error surface

on each iteration

- current weights define a point in this space
- find direction in which error surface descends most steeply
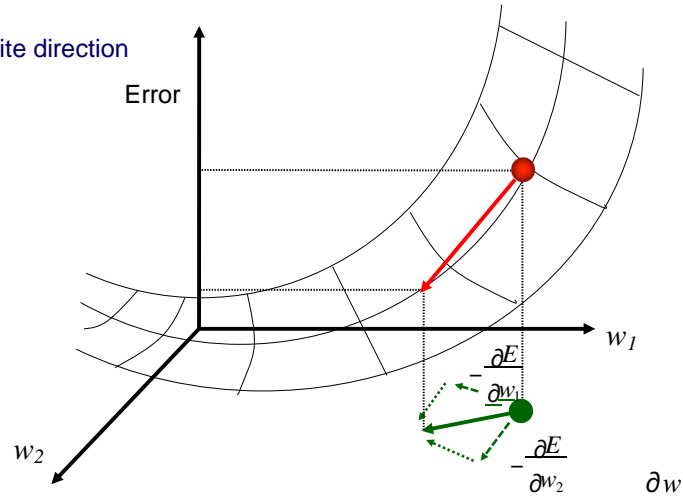- take a step (i.e. update weights) in that direction

# Gradient descent in weight space

calculate the gradient of $E$: $\quad \nabla E(w) = \left[\dfrac{\partial E}{\partial w_0}, \dfrac{\partial E}{\partial w_1}, \square, \dfrac{\partial E}{\partial w_n}\right]$
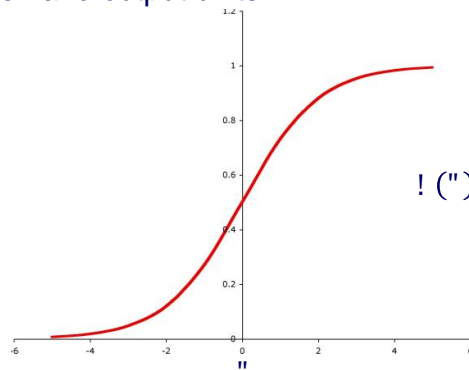
take a step in the opposite direction

$$\Delta w = -\eta \; \nabla E(w)$$

$$\Delta w_i = -\eta \; \dfrac{\partial E}{\partial w_i}$$
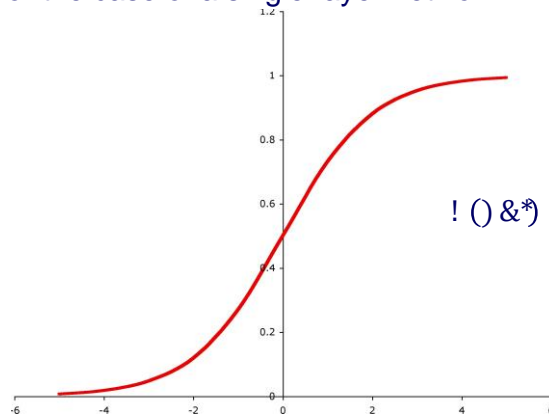


# The sigmoid function

- to be able to differentiate $E$ with respect to $w_i$, our network must represent a continuous function
- to do this, we can use *sigmoid functions* instead of threshold functions in our hidden and output units



$$\sigma(z) = \dfrac{1}{1 + e^{-z}}$$

# The sigmoid function

for the case of a single-layer network

$$! () \&^*) = \frac{1}{1 + \&^{'\,(+,\,-\,\Sigma_/ + /1_/)}}$$

$$) \&^* = 2_3 + 4 \quad 2_5 6_5$$

# BatchNN training

**given**: network structure and a training set $D = \{(x_{(1)}, y_{(1)}) \square (x_{(m)}, y_{(m)})\}$

initialize all weights in $w$ to small random numbers

until stopping criteria met do initialize the error

$$E(w) = 0$$

for each $(x^{(d)}, y^{(d)})$ in the training set input $x^{(d)}$ to the network and compute output $o^{(d)}$ increment the error

$$E(w) = E(w) + \frac{1}{2}\left(y^{(d)} - o^{(d)}\right)^2 \text{ calculate the gradient}$$

$$\nabla E(w) = \left[\frac{\#\ \partial E}{\partial}, \frac{\partial E}{\partial}, \frac{\partial E}{\partial w_n}\right]_{w_0\, w_1 \square,}$$

update the weights

$$\Delta w = -\eta\ \nabla E(w)$$

# Online vs. Batch learning

- Standard gradient descent (batch training): calculates error gradient for the <u>entire training set</u>, before taking a step in weight space

- *Stochastic gradient descent* (online training): calculates error gradient for a single instance (or a small set of instances, a "mini batch"), then takes a step in weight space
  - much faster convergence
  - less susceptible to local minima

# OnlineNN training (Stochastic Gradient Descent)

**given**: network structure and a training set $D=\{(x_{(1)}, y_{(1)})\ldots(x_{(m)}, y_{(m)})\}$

initialize all weights in $w$ to small random numbers until

stopping criteria met do

for each $(x^{(d)}, y^{(d)})$ in the training set input $x^{(d)}$ to the network and compute output $o^{(d)}$ calculate the

error $E(w) = \dfrac{1}{2}(y_{(d)} - o_{(d)})_2$ calculate the gradient

$$\nabla E(w) = \% \underline{\quad} \underline{\quad} \begin{bmatrix} \dfrac{\partial E}{\partial w_0}, \dfrac{\partial E}{\partial w_1}, \square, \dfrac{\partial E}{\partial w_n} \end{bmatrix}$$

update the weights
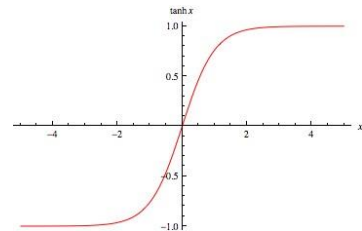
$$\Delta w = -\eta \; \nabla E(w)$$

# Other activation functions

• the sigmoid is just one choice for an *activation function* •

there are others we can use including

hyperbolic tangent

$$\varphi(6) = \tanh 6 = \frac{2}{1 + e^{-2x}} - 1$$

rectified linear (ReLU)

$$\varphi(6) => \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$$

# Other objective functions

• squared error is just one choice for an *objective function* •

there are others we can use including cross entropy

$$E_F = \sum_{H \in N} -y^{(d)} \ln\big(o^{(d)}\big) - \big(1 - G^H\big) \ln\big(1 - J^{(H)}\big)$$
multiclass

cross entropy

$$E_F(\ ) = -\sum_{H \in N} \sum_{5VW}^{\# \ QRSTTUT} y_i^{(d)} {}_{O)} \big(J_{5(H)}\big)$$

# Convergence of gradient descent

- gradient descent will converge to a minimum in the error function
- for a <u>multi-layer network</u>, this may be a *local minimum* (i.e. there may be a "better" solution elsewhere in weight space)



- for a <u>single-layer network</u>, this will be a global minimum (i.e. gradient descent will find the "best" solution)
- Recent analysis suggests that local minima are probably rare in high dimensions; saddle points are more of a challenge      [Dauphin et al., NIPS 2014]